# 1 Review

## 1.1 Online Learning

Recall the goal of online learning (Figure 1) is to make a series of accurate predictions given a sequence of inputs. At each time step $t$, the learner takes in the input and any additional information (such as expert predictions). It then makes a prediction $\hat{y}$ that validates given a true label $y$. Based on the calculated loss between these two, it updates its parameters for the next time step.
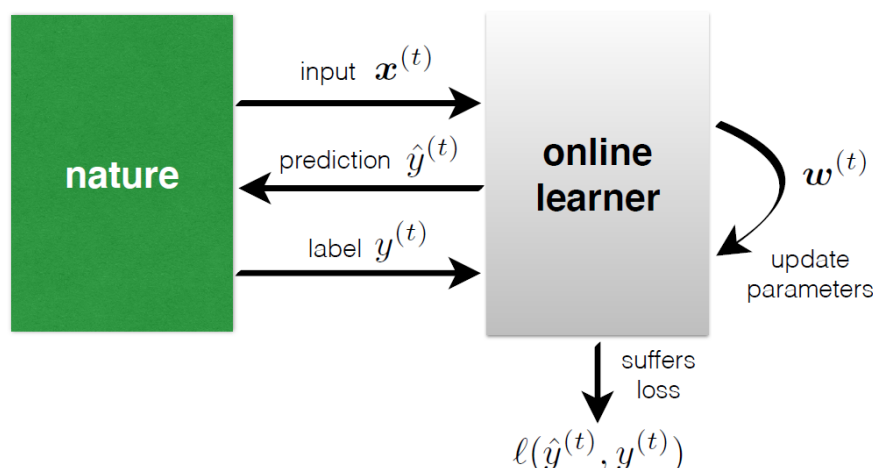


Figure 1: Online Learning (Lecture 6, Slide 16)

Some topics we have discussed previously are the mistake and regret bounds of the online learning algorithms. Mistake bounds are the upper limit on the expected number of mistakes the algorithm can make. Regret bounds compare the expected loss of the learner with the best hypothesis $h \in \mathcal{H}$.

We have seen online learning algorithms in two different categories thus far. *Greedy (Consistent)*, *Halving*, *Weighted Majority*, and *Randomized Weighted Majority* Algorithms all fall under **Prediction with Expert Advice**. In this category, the learner is given a set of expert predictions at each time step. The second category is **Online Linear Classification**, and we have seen the *Online Perceptron* and *Winnow* Algorithm in this category.

## 1.2 Convexity

**Definition 1.** A set $S$ is **convex** if for all $w, v \in S$, $\alpha w + (1 - \alpha)v \in S$ for all $\alpha \in [0, 1]$ In two

dimensions, this means that a set $S$ is convex if a line connecting any two points is contained within the set. Figure 2 is an example of a convex set, and for any two points $X$ and $Y$, the line connecting it is always contained within the set.
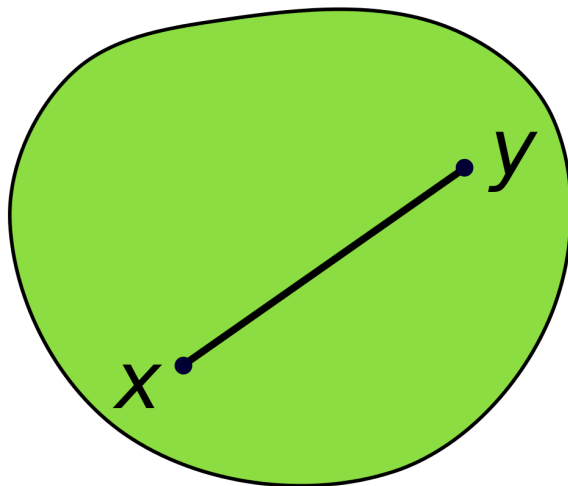


Figure 2: Convex Set (image from Wikipedia)

**Definition 2.** A function $f : S \in \mathbb{R}$ is **convex** if for all $w, v \in S$, $f(\alpha w + (1 - \alpha)v) \leq \alpha f(w) + (1 - \alpha)f(v)$ for all $\alpha \in [0, 1]$

In two dimensions, this means that a convex function $f$ between any two points is upper-bounded by a line connecting the two points. Or that a line connecting two points on a convex function always lies above the function. Figure 3 is an example of a convex function. The line connecting any two points $x_1$ and $x_2$ always lies above the function.
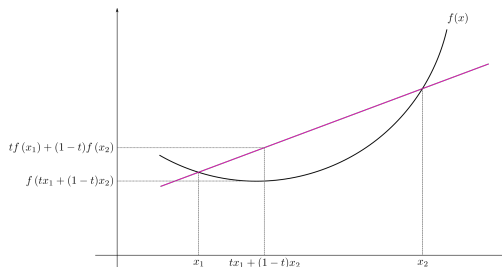


Figure 3: Convex Function (image from Wikipedia)

## 2   Summary

### 2.1   Optimization

Optimization involves minimizing or maximizing a function by choosing valid inputs and computing the value of the function on those inputs.

**Definition 3.** A (constrained) optimization problem can be formulated as $\min f(x)$, with constraints $g(x) \leq 0;\ h(x) = 0\ x \in S$

$g(x)$ is the inequality constraint, $h(x)$ is the equality constraint, and $S$ is the optimization domain. There are three solution methods for solving optimization problems, each with benefits and drawbacks.

## Analytic Solution

The analytic solution such as to $\nabla f(x) = 0$ can provide a global and fast solution to the optimization problem. However, it can be very difficult and sometimes impossible to compute.

## Brute Force Search

Brute force search involves testing all possible valid inputs, exhaustively, to find the optimal value. This approach has the benefit of providing a global solution, but is usually very slow for a large input space and is sometimes impossible depending on the formulation of the problem.

## Numerical Methods

Numerical methods, such as gradient descent, are a popular and sometimes fast way to solve optimization problems. However, these methods usually do not have guarantees for finding a global solution.

Numerical methods can be used to solve online convex optimization by:

1. Receiving one data sample

2. Observing the loss

3. Update the parameters (weights)

This process, using an online learner, is illustrated in Figure 4.
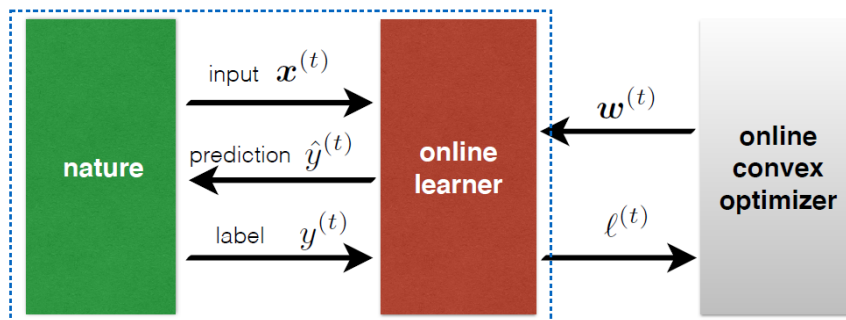


Figure 4: Online Convex Optimization (Lecture 6, Slide 18)

3

## 2.2 Lipschitz Continuity

A function $f$ is called **L-Lipschitz** over a set $S$ with respect to a norm, if the following holds for all $\boldsymbol{u}, \boldsymbol{w} \in S$

$$|f(\boldsymbol{u}) - f(\boldsymbol{w})| \leq L||\boldsymbol{u} - \boldsymbol{w}||$$

From Wikipedia [3]: Intuitively, Lipschitz continuity is a measure for how fast a function can change: There exists a real number such that, for every pair of points on the graph of this function, the absolute value of the slope of the line connecting them is not greater than this real number.
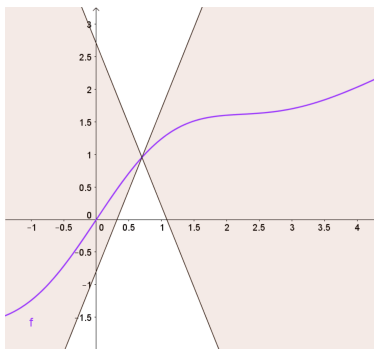


Figure 5: Lipschitz continuity visualization. From [3].

If a function is Lipschitz continuous, there exists a double cone (depicted in white in 5) whose origin can be moved along the graph such that the whole graph always stays outside of the double cone.

## 2.3 Convexification

There exist many problems and learning applications in which the loss function is not convex. We can however apply a few tricks to turn a non-convex loss function into a convex one. This process is called **convexification**, and we will summarize two methods in the following sections: Convexification by randomization and convexification by surrogate loss.

**Convexification by randomization**

We will now consider the Weighted Majority Algorithm (WMA) (see Lecture 3 for more details). This algorithm predicts based on the majority vote

$$\hat{y} = \text{sign}\langle \boldsymbol{x}^{(t)}, \boldsymbol{w}^{(t-1)}\rangle$$

and we receive a loss

$$l^{(t)} = \mathbf{1}[\hat{y}^{(t)} \neq y^{(t)}]$$

This algorithm is a bounded-regret algorithm. If we examine the parameter space, we realize that the parameter space is the set of positive real numbers, $\mathbb{R}^+$. This is a convex set. However, the loss function is either zero or one, which is a non-convex function, as shown in Fig. 6. We will now
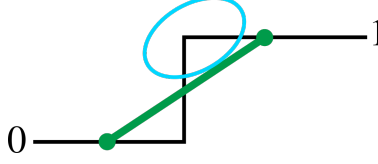
Figure 6: Zero-One loss. Since the circled portion of the function is above the line segment, this problem is a non-convex function.

show that we can use randomization to make this loss function convex. If we look at the prediction rule of the Randomized Weighted Majority Algorithm (RWMA)

$$h \sim \text{MULTINOMIAL}(\boldsymbol{w}^{(t)}/\Phi^{(t)})$$

$$\hat{y}^{(t)} = h(\boldsymbol{x}^{(t)})$$

we predict based on a multinomial distribution over the relative expert weights $w^{(t)}$. The parameter space here is the set of positive real numbers, $\mathbb{R}^+$. The loss function for RWMA is given by

$$l^{(t)} = E_p[\mathbf{1}[y^{(t)} \neq \hat{y}_n^{(t)}]] = \sum_{p_n} p_n \mathbf{1}[y^{(t)} \neq \hat{y}_n^{(t)}])$$

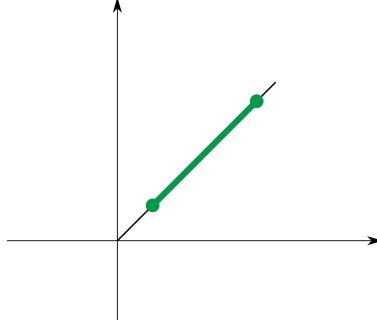which can be visualized as shown in Fig. 7.



Figure 7: Linear loss, a convex function.

We observe that this is a linear function, which is convex. We have therefore shown that by introducing randomization to the algorithm, we have transformed the loss function into a convex function. It should be noted that randomization also turned this algorithm from being bounded-regret into a no-regret algorithm.

**Convexification by surrogate loss**

A further method to transform a non-convex problem into a convex problem is to replace the original non-convex loss function with a **surrogate** loss function.

In order to apply this method, there are two requirements:

- The surrogate loss upper bounds the original loss

- The surrogate loss is convex

If we now minimize the surrogate loss, we are minimizing an upper bound of the original loss.

For WMA, we can introduce the *hinge* loss as a surrogate loss function. The hinge loss is given by the following equations.

$$
\begin{aligned}
l^{(t)} &= \mathbf{1}[\hat{y}^{(t)} \neq y^{(t)}] \\
&= \mathbf{1}[1 - y^{(t)}\langle \boldsymbol{w}^{(t-1)}, \boldsymbol{x}^{(t)}\rangle > 1] \\
&\leq \max\big[0, 1 - y^{(t)}\langle \boldsymbol{w}^{(t-1)}, \boldsymbol{x}^{(t)}\rangle > 1\big] \\
\tilde{l}^{(t)} &= \max\big[0, 1 - y^{(t)}\langle \boldsymbol{w}^{(t-1)}, \boldsymbol{x}^{(t)}\rangle > 1\big]
\end{aligned}
$$

As shown in Fig. 8, the hinge loss $\tilde{l}^{(t)}$ is a convex function. We have therefore shown that the surrogate loss upper bounds the original loss, and turned this problem into a convex problem.
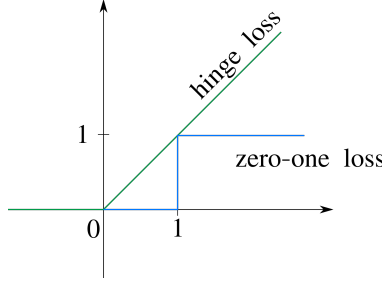


Figure 8: Hinge loss upper bounds zero-one loss and is a convex function.

## 2.4   Follow the Leader

We will now discuss an algorithm for online convex optimization problems. Follow the leader (FTL), also termed *Fictitious Play*, was first introduced by G. W. Brown in 1951 [2, 1].

---
**Algorithm 1** Follow the leader
---
1: **function** FOLLOW THE LEADER
2:     **for** $t = 1, 2 \cdots, T$ **do**
3:         $\boldsymbol{w}^{(t)} = \arg\min_{\boldsymbol{w} \in W} \sum_{i=1}^{t-1} f^{(i)}(\boldsymbol{w})$
4:         RECEIVE$(f^{(T)} : W \to \mathbb{R})$
---

The main idea of this algorithm is that the learner should go by the best choice seen so far. We can plug the weight update rule (line 3 in Algorithm 1) into the online convex optimization as shown in Fig. 4.

We will now derive the regret bound. For this purpose, we will introduce the concept of the **one-step look ahead cheater**.

6

The cumulative regret is given by

$$R(\boldsymbol{u}) = \sum_t [f^{(t)}(\boldsymbol{w}^{(t)}) - f^{(t)}(\boldsymbol{u})]$$

where $\boldsymbol{u}$ is any expert, $\boldsymbol{w}^{(t)}$ is the set of parameters at time $t$, $f^{(t)}$ the loss function returned to the online learner by nature. It is difficult to reason about $\boldsymbol{u}$ as there are infinite possibilities. We can however upper bound the regret by introducing a *one step look ahead cheater* $\boldsymbol{w}^{(t+1)}$:

$$R(\boldsymbol{u}) \leq \sum_t [f^{(t)}(\boldsymbol{w}^{(t)}) - f^{(t)}(\boldsymbol{w}^{(t+1)})]$$

We will now prove that this inequality holds true. We start by subtracting $f^{(t)}(\boldsymbol{w}^{(t)})$ from both sides of the inequality, which leaves us with

$$\sum_{t=1}^{T} f^{(t)}(\boldsymbol{w}^{(t+1)}) \leq \sum_{t=1}^{T} f^{(t)}(\boldsymbol{u})$$

This means that the loss of a series of one step cheaters is less than or equal to the loss of any other single parameter $\boldsymbol{u}$.

We can prove that this inequality holds true by Induction: We assume that the inequality holds for $T-1$

$$\sum_{t=1}^{T-1} f^{(t)}(\boldsymbol{w}^{(t+1)}) \leq \sum_{t=1}^{T-1} f^{(t)}(\boldsymbol{u}) \quad \forall \boldsymbol{u}$$

If we add $f^{(T)}(\boldsymbol{w}^{(T+1)})$ to both sides, we obtain

$$\sum_{t=1}^{T} f^{(t)}(\boldsymbol{w}^{(t+1)}) \leq f^{(T)}(\boldsymbol{w}^{(T+1)}) + \sum_{t=1}^{T-1} f^{(t)}(\boldsymbol{u})$$

Note that the term can be consumed within the summation on the left hand side, but not on the right hand side, since the parameters are different. But since we assume that the inequality is true for all $\boldsymbol{u}$, we simply set it to $\boldsymbol{u} = \boldsymbol{w}^{(T+1)}$, and can write

$$\sum_{t=1}^{T} f^{(t)}(\boldsymbol{w}^{(t+1)}) \leq f^{(T)}(\boldsymbol{w}^{(T+1)}) + \sum_{t=1}^{T-1} f^{(t)}(\boldsymbol{w}^{(T+1)})$$

$$\leq \sum_{t=1}^{T} f^{(t)}(\boldsymbol{w}^{(T+1)})$$

From this inequality we can derive that the loss of using the one-step look ahead cheater at each time step is upper bounded by the loss of using a single one-step look ahead cheater at the end of the sequence for all time steps. Given the definition of the minimizer, $\boldsymbol{w}^{(T+1)} = \arg\min_{\boldsymbol{u} \in S} \sum_{t=1}^{T} f^{(t)}(\boldsymbol{u})$

this means we can upper bound the right hand side of the inequality by all $\boldsymbol{u}$:

$$\sum_{t=1}^{T} f^{(t)}(\boldsymbol{w}^{(t+1)}) \leq \sum_{t=1}^{T} f^{(t)}(\boldsymbol{u})$$

We can now plug this back into our original upper bound on the regret, and finally obtain

$$\begin{aligned} R(\boldsymbol{u}) &= \sum_{t} [f^{(t)}(\boldsymbol{w}^{(t)}) - f^{(t)}(\boldsymbol{u})] \\ &\leq \sum_{t} [f^{(t)}(\boldsymbol{w}^{(t)}) - f^{(t)}(\boldsymbol{w}^{(t+1)})] \end{aligned}$$

which concludes our proof.

# References

[1] U. Berger. Brown's original fictitious play. *Journal of Economic Theory*, 135(1):572–578, 2007.

[2] G. W. Brown. Iterative solution of games by fictitious play. *Activity analysis of production and allocation*, 13(1):374–376, 1951.

[3] W. contributors. Lipschitz continuity — Wikipedia, the free encyclopedia, 2020. [Online; accessed 18-February-2021].